

How? → Algorithm → Set of Rules → Cleaned Data

Agenda

→ {Elliptical Envelope}

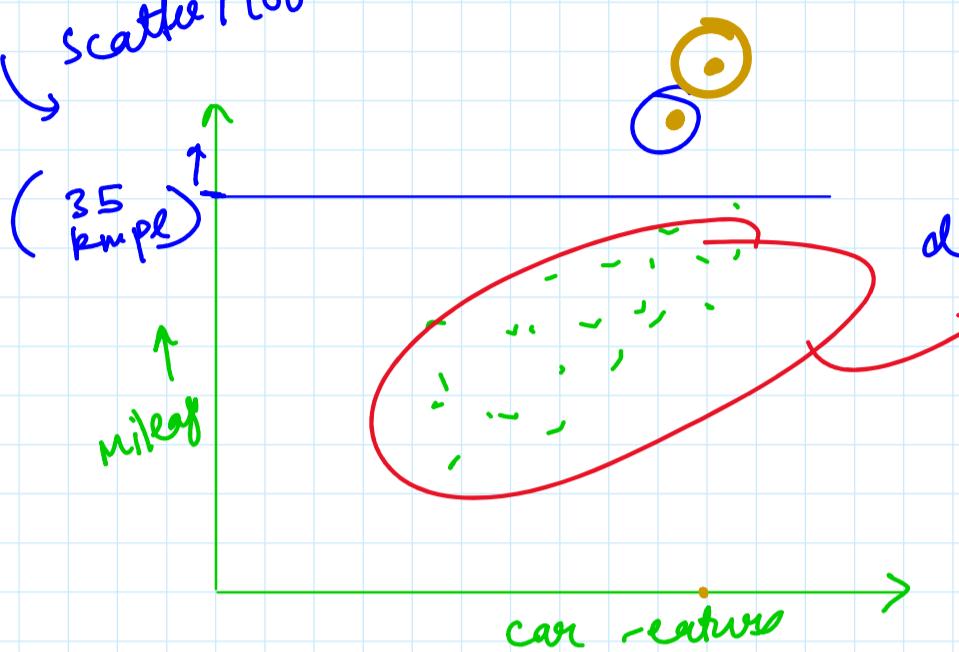
→ Isolation Forest

→ LOF (Local outliers factor) (Drive)

Simple Methods

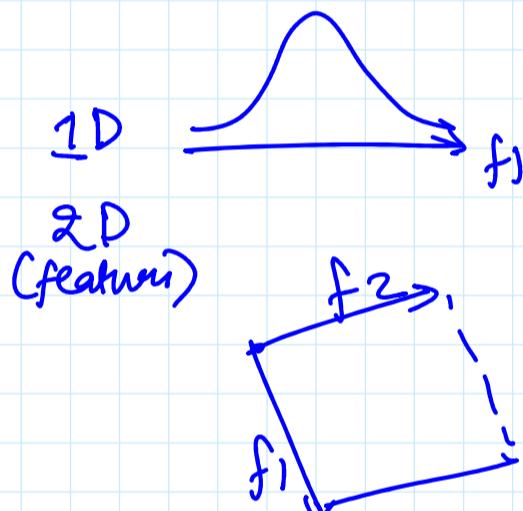
- Box Plot
- Z score $\rightarrow (Z = \frac{x-\mu}{\sigma}) \rightarrow 3.2 [Z = 3]$
- Visualise (Remove manually)
 - Scatter Plot → Filter Data

Scatter Plot

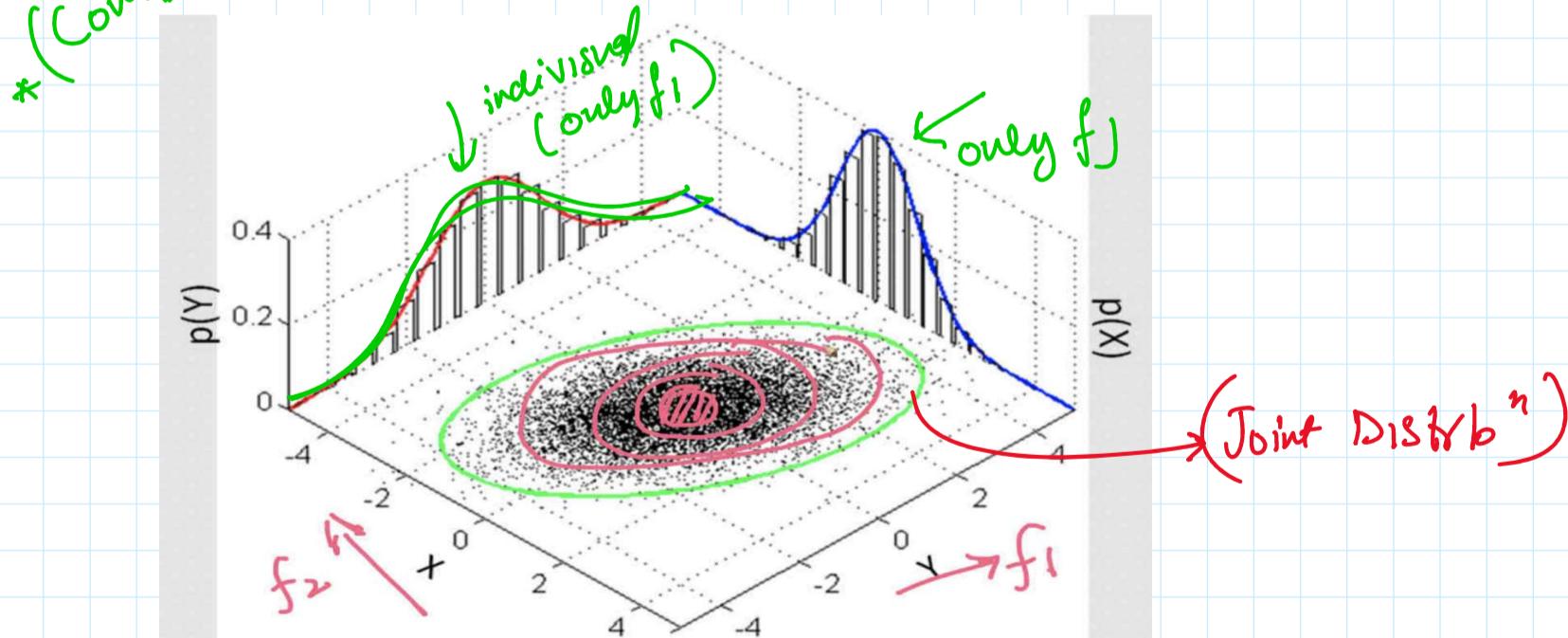


$$df = df[df['mileage'] \leq 35] \quad T/F$$

Elliptical Envelope



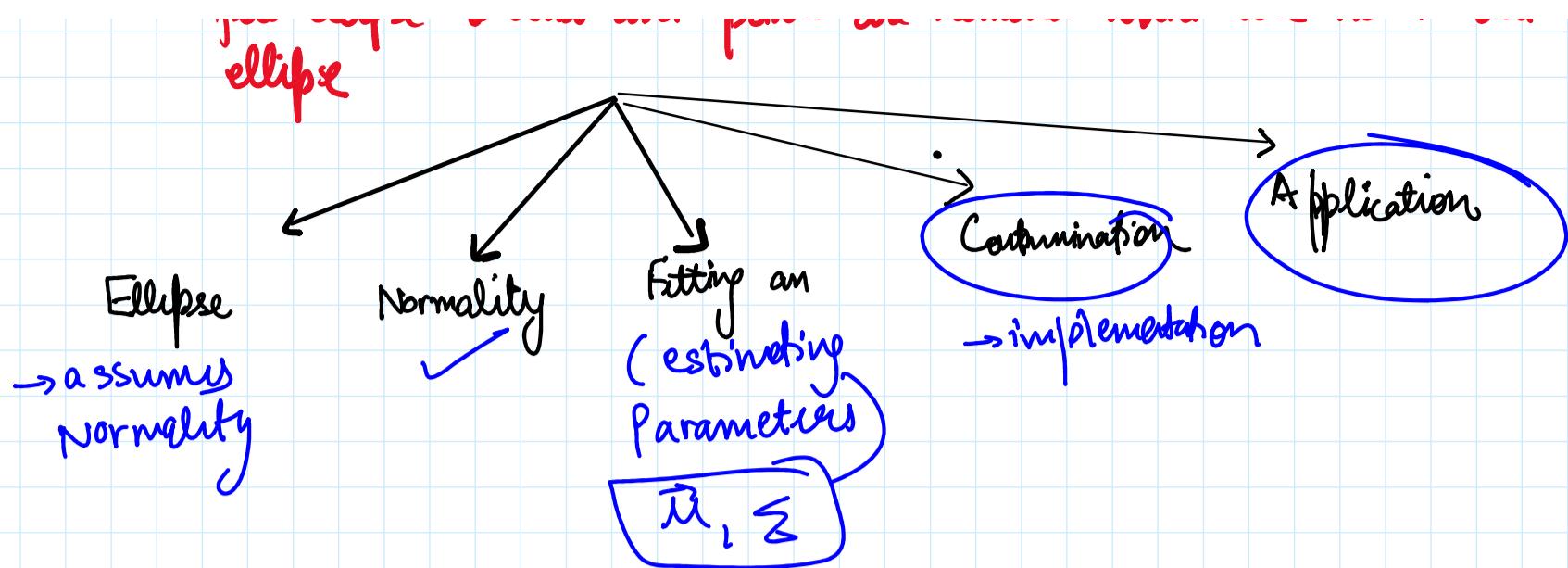
* (Combined)



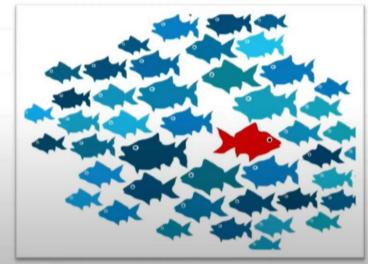
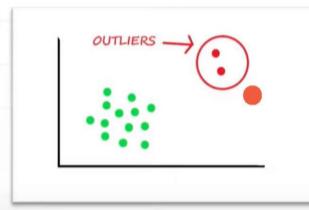
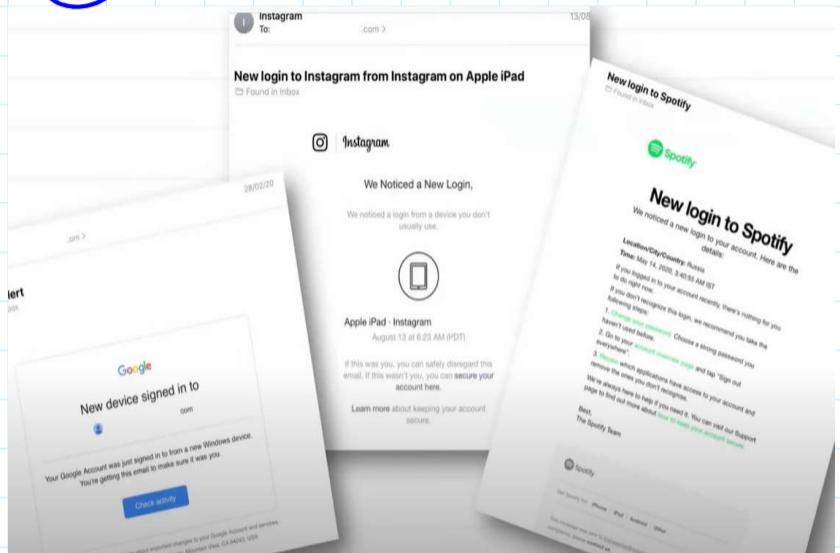
assume
(Data \rightarrow Gaussian)

Elliptical Envelope \rightarrow statistical / Multivariate / Outlier Detection / Gaussian

fits ellipse to data and points are removed which are not in that ellipse



(X) Blocks that request



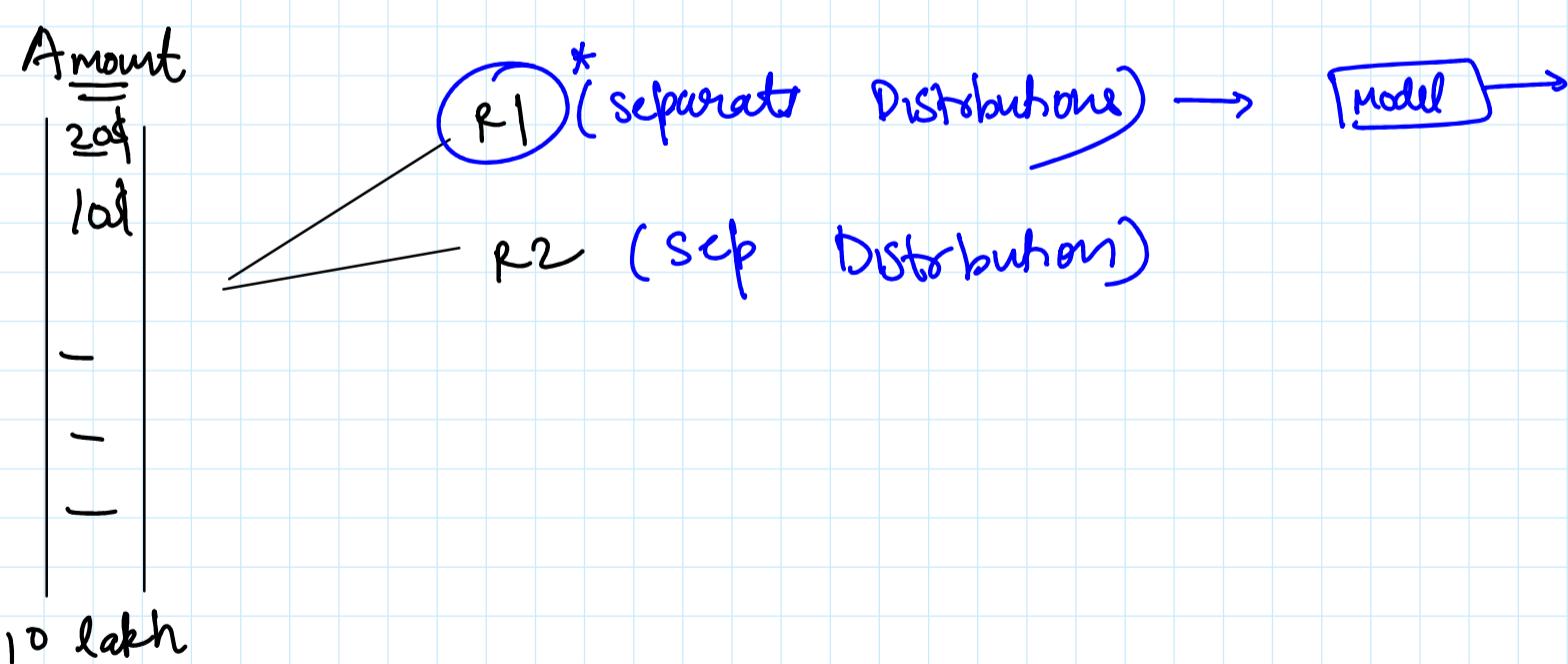
threshold = $Z \rightarrow 3$ ($Z \geq 3$) (outlier)

Payments → (Red money) $Z \rightarrow < 3$ (safe)

> 3

$= 3$

(you don't know)

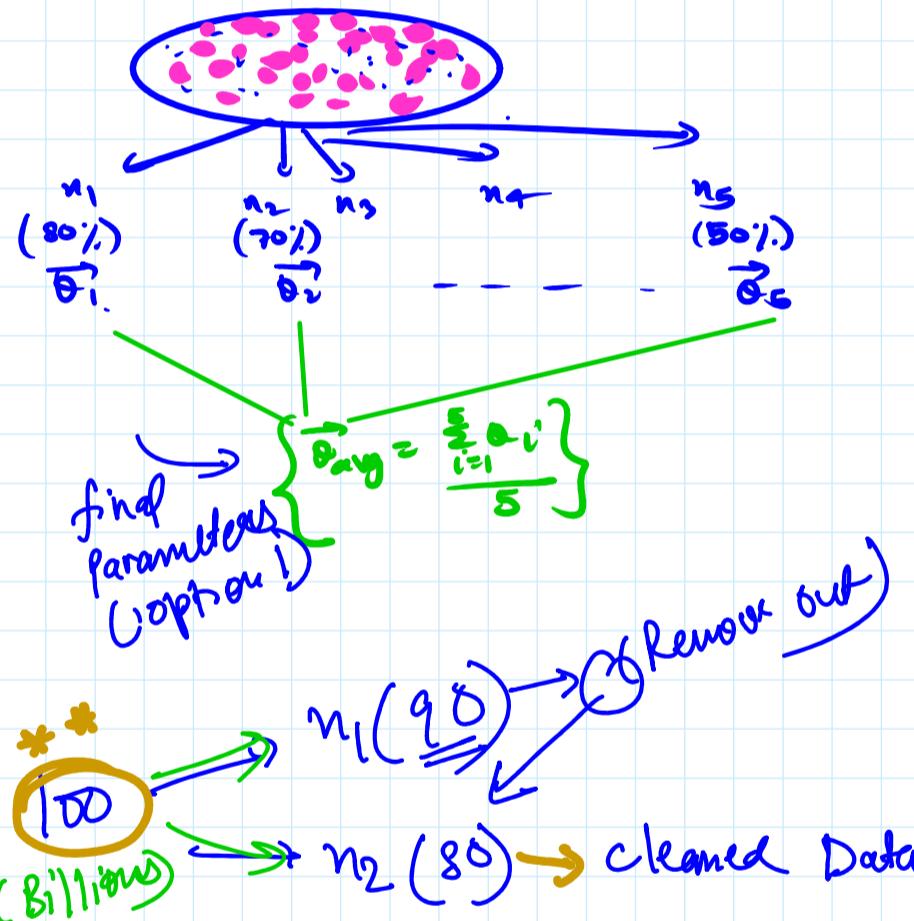
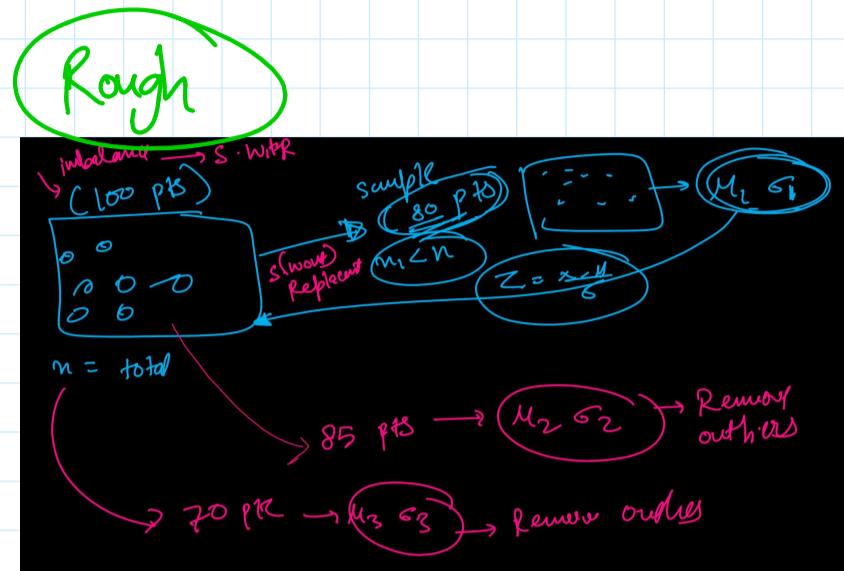
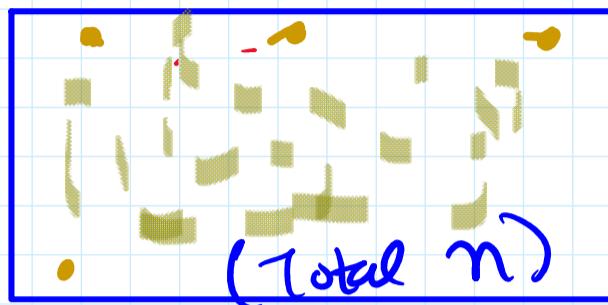


Fit a Distribution (one distribution) → Dimension is high
→ (Data) → fit distrib'n
(we have)

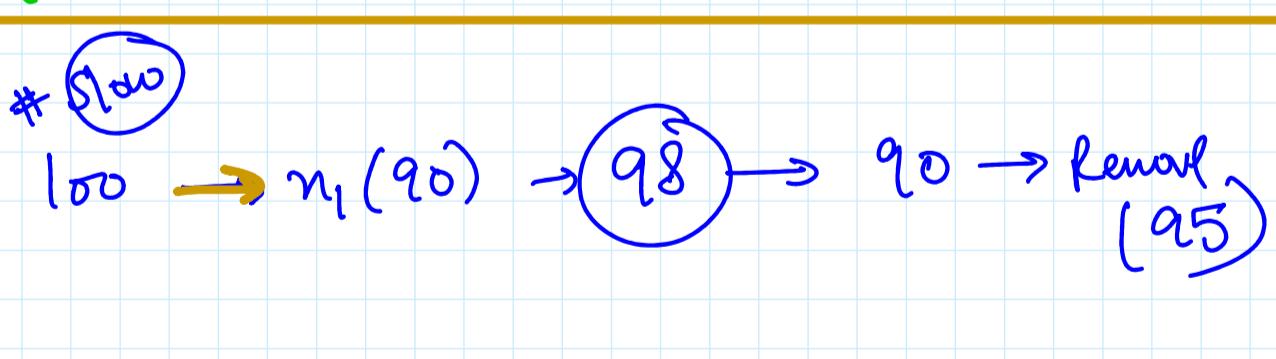
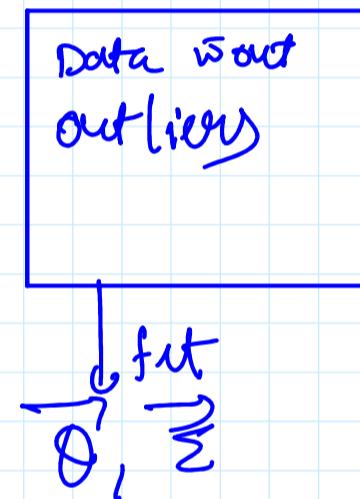
New Point → Z score → threshold → Anomaly

* RANSAC

RANSAC: RANSAC is an iterative method used primarily in computer vision and machine learning to estimate parameters of a mathematical model from a set of observed data that contains outliers. It's particularly useful when you expect a significant fraction of the data to be outliers.



option (II)



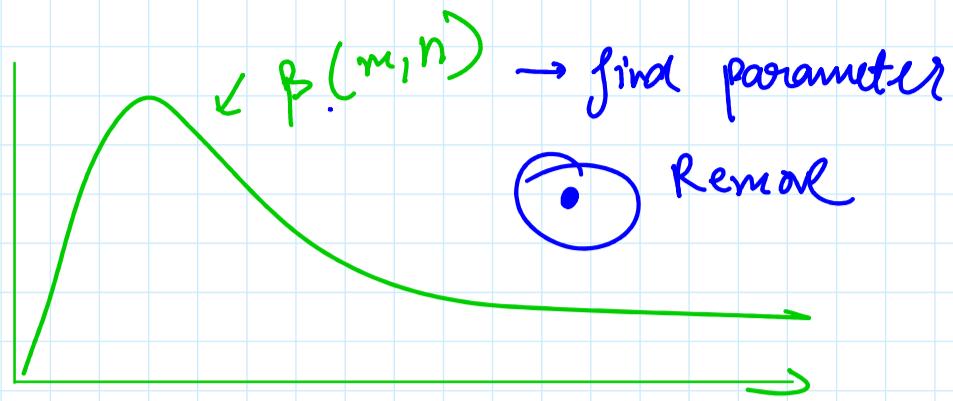
E-E assume Gaussian \rightarrow convert $X \rightarrow \log(X) = Y$

$$Y \sim N(\mu, \sigma)$$

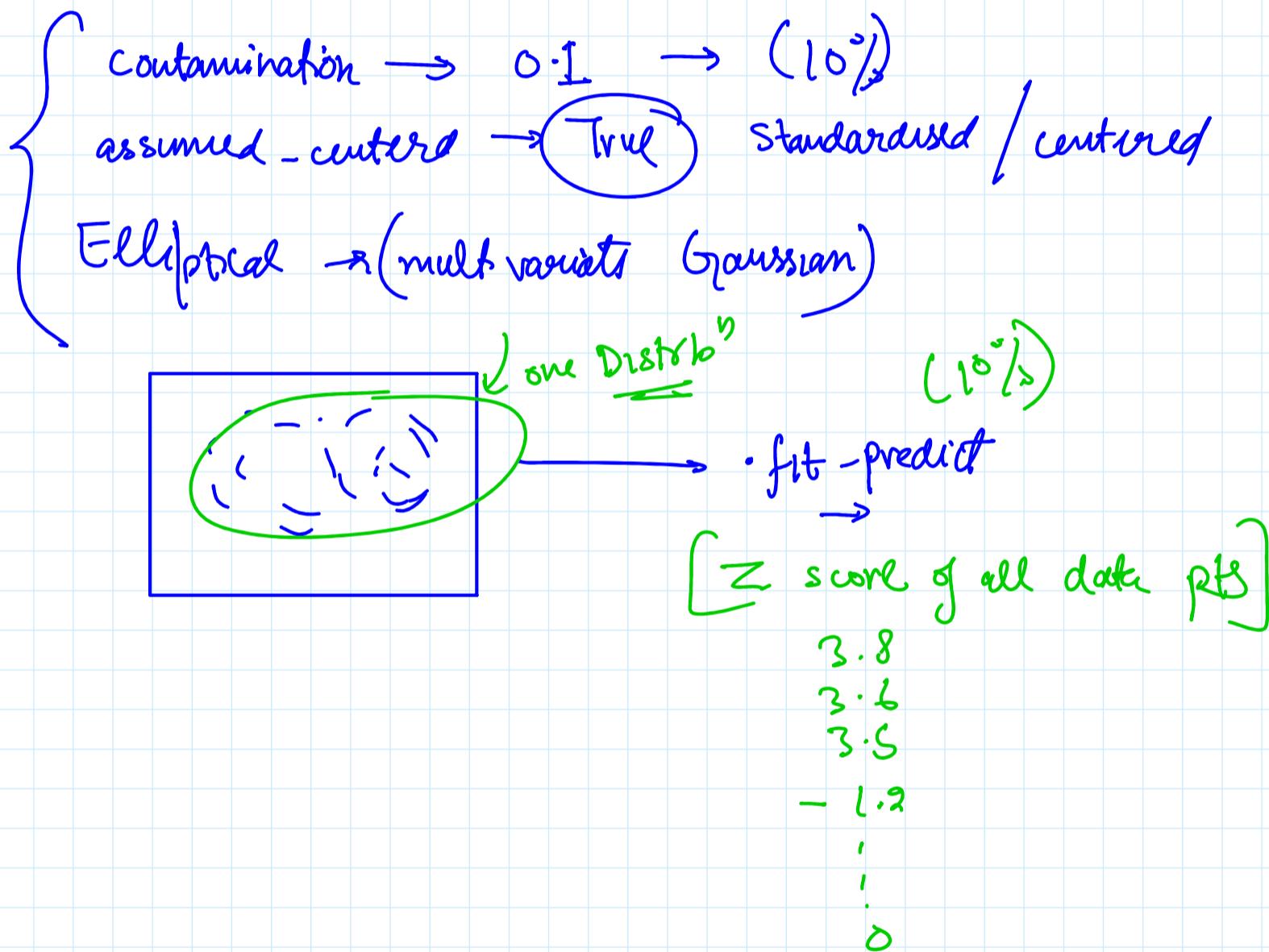
Q. For non-gaussian, do we need to convert into gaussian distribution?

- While the elliptical envelope method makes an assumption that the distribution is gaussian, the strategy can be applied to any distribution
- As long as we know any distribution and its parameters θ , we can extend our strategy to use RANSAC and estimating parameters θ
- The can be other distributions such as multivariate poissons, multivariate log normals, etc., but we don't use them as much
- One other strategy is to convert them into gaussians but we don't necessarily have to.

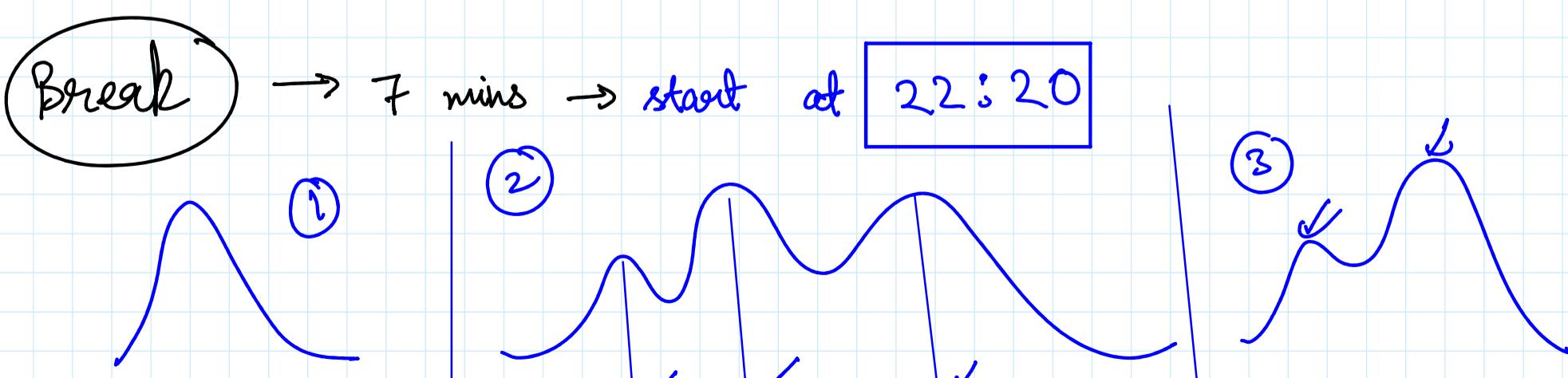
- The can be other distributions such as multivariate poissons, multivariate log normals, etc., but we don't use them as much
- One other strategy is to convert them into gaussians but we don't necessarily have to.
- As long as there is any distribution we can calculate the probability of $x_i \in X$ using PMF/PDF

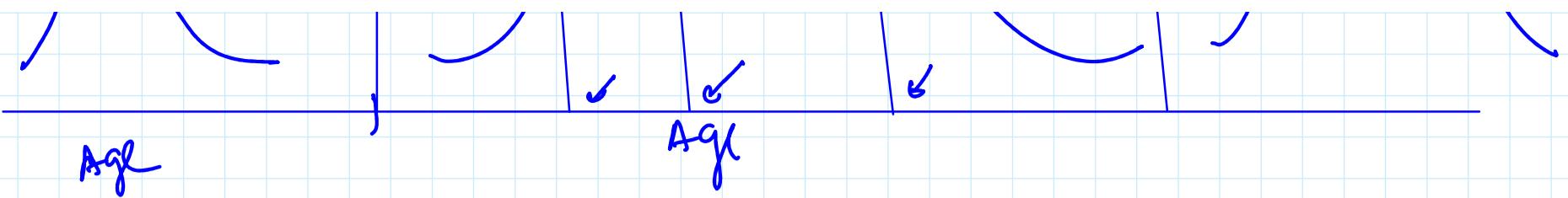


$\text{sk-Learn} \rightarrow$ (Package) $\cdot \text{fit}$ $\cdot \text{predict}$
 labels the outliers



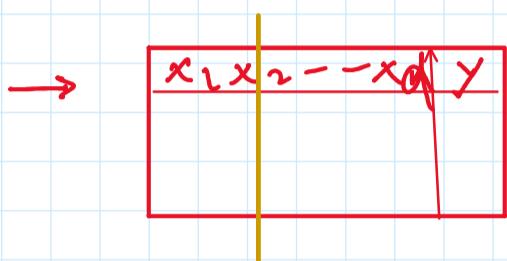
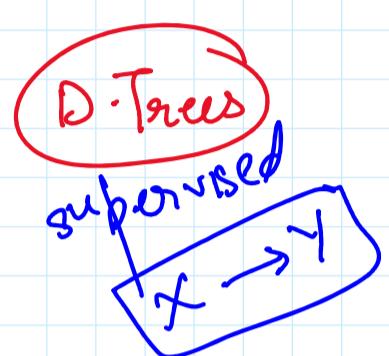
Disadvantages \rightarrow Disadvantages
 You might have find the concept of elliptic envelope quite straightforward, but there are limitations.
 • It cannot be used non-unimodal data
 • It is specifically for multivariate gaussians
 If the data fails to meet the assumptions of unimodal and multivariate gaussian, the whole things crashes



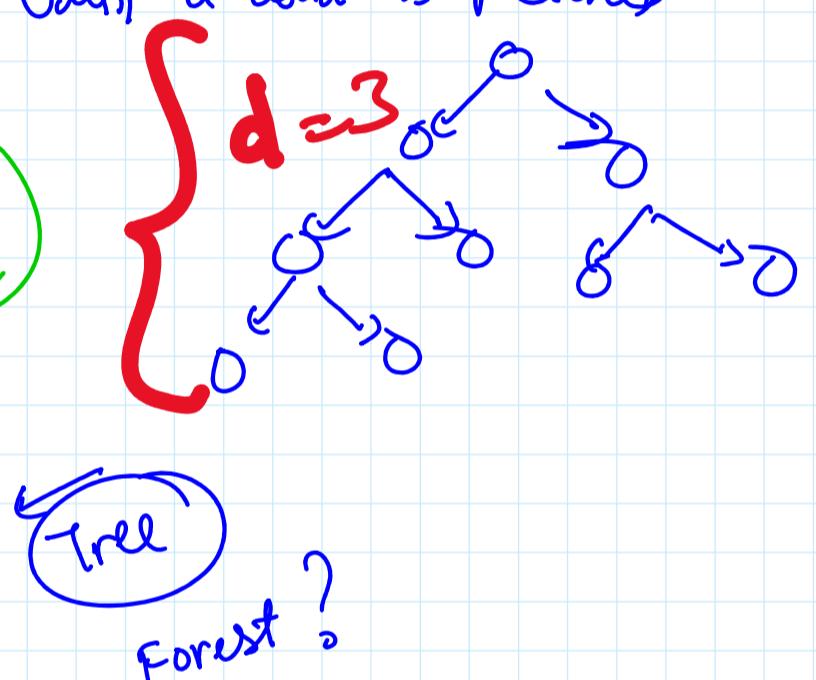
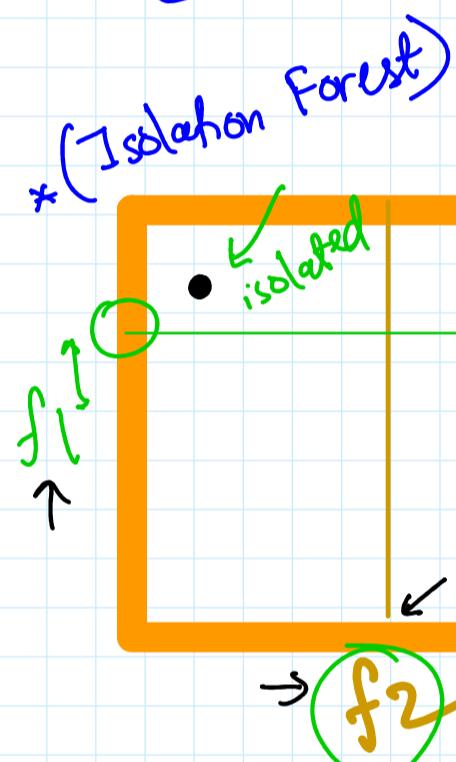


Isolation Forest (Categorical + Numerical)

Concept → Behaviour $\begin{cases} \rightarrow \text{outliers} \\ \rightarrow \text{Inliers} \end{cases}$ is different



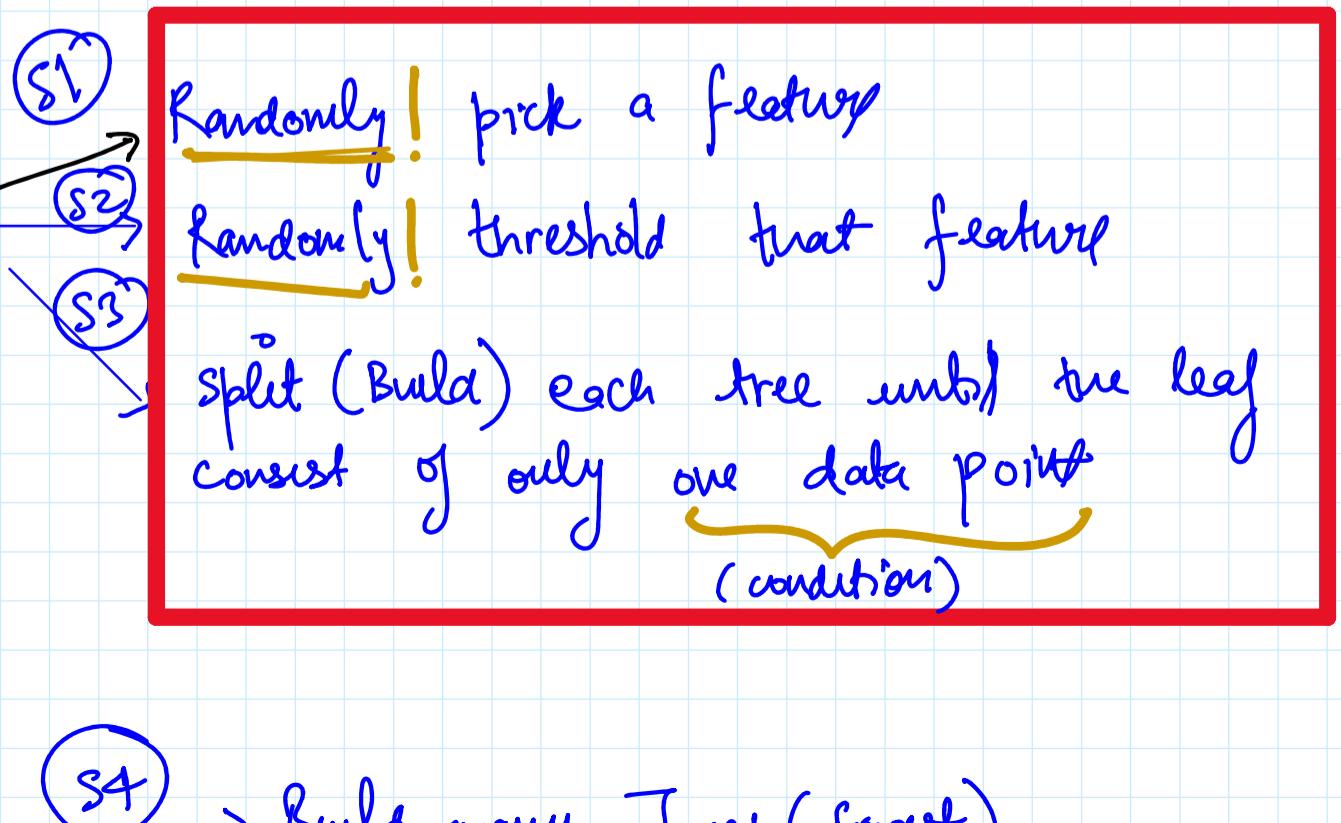
"best"
Choose feature
→ Choose a split threshold
→ Keep splitting repeatedly
→ Until a "cond" is reached



Logic → Outliers are isolated easily than inliers

Unsupervised → NO Y (Target)

Isolation Forest Algo

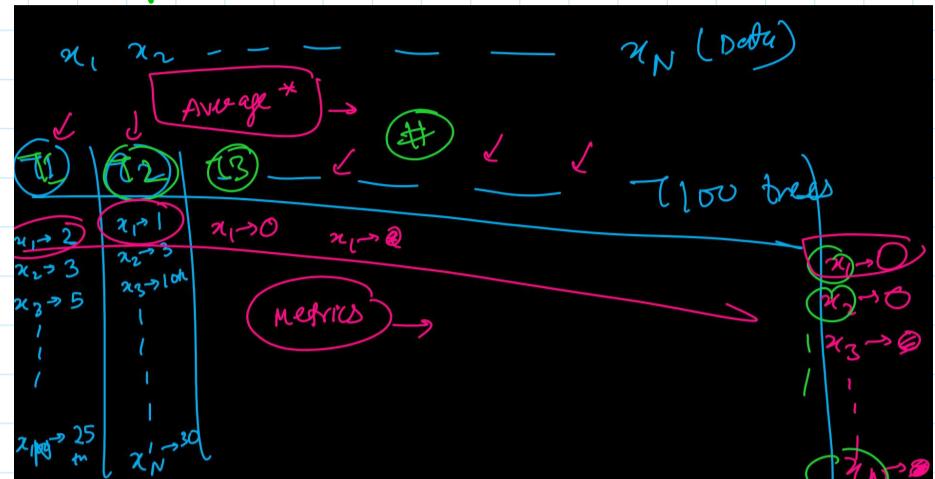


(S4) → Build many Trees (forest)

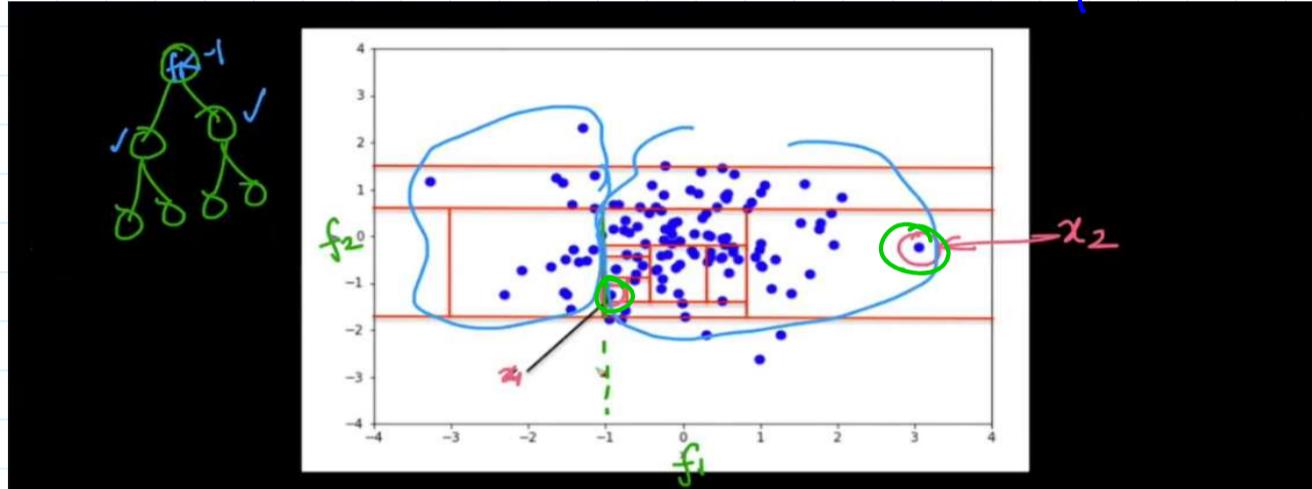
Rough (2)



Rough (3)



Axis Parallel Split Problems



. What is the whole idea behind Isolation Forest?

So, to sum it up, the idea behind Isolation forest is,

- On an average outliers have lower depth in the random trees
- On an average, inliers have more depth in the random trees

How can we evaluate Isolation Forests?

- Imagine, we have build 100 random trees. For each point x_i in the dataset, we can get an average depth.
- We use this average depth to convert into a metric.
- Apart from this, there are lot of different metrics, that people have came up with over the years
- But, the basic intuition is that lesser the average depth, higher likelihood is there that it is an outlier

* (sample some rows)

But, what if the number of datapoints is large?
Wouldn't it mess up the Isolation Forests?

- iForests can be made on subset of samples.
- We use this subset as train dataset and the rest of the data as test database

But, how do we decide average depth for a point to be classified as an outlier?

- There is no one metric specifically used for average depths in iForests. At the end, whichever metric you use, it is based on the threshold.
- There are a lot of metrics that researchers have came up with over the years.
- But, studying them in this lecture is out of scope.

learn

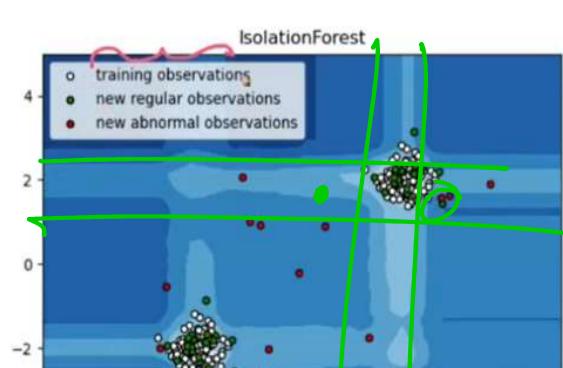
Prev Up Next

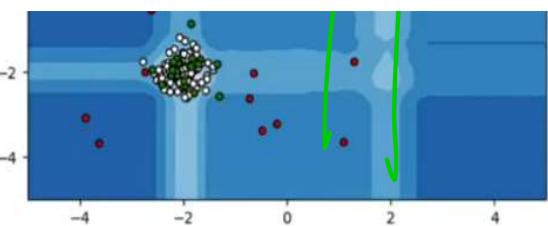
scikit-learn 1.1.1
Other versions

Please cite us if you use the software.

IsolationForest example

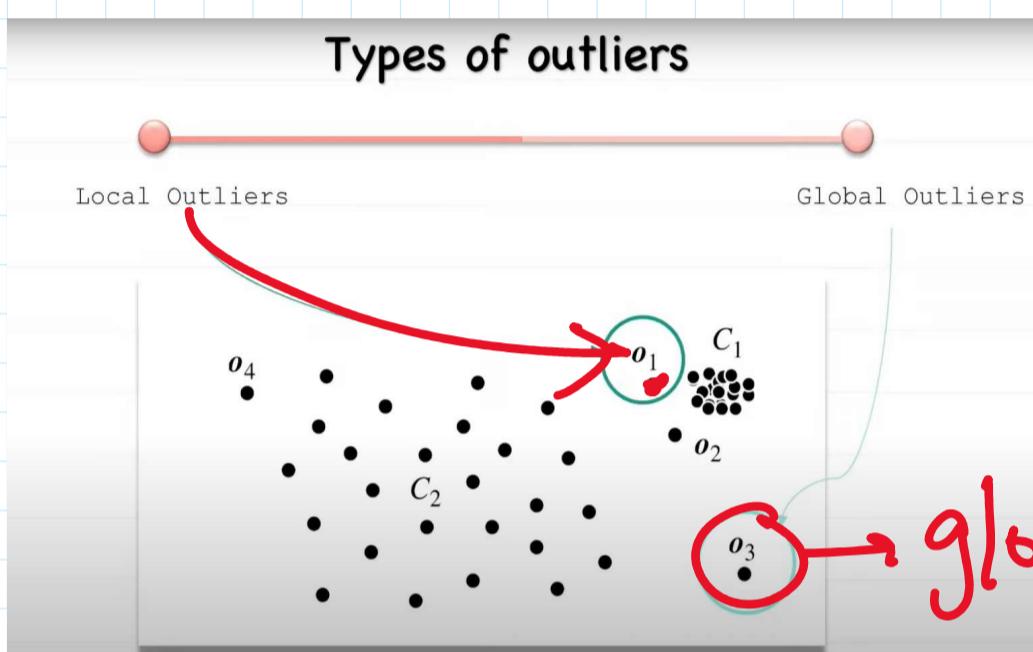
Random partitioning produces noticeable shorter paths for anomalies. Hence, when a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies.





D.1
 ④ Axis Parallel Splits → Cannot capture non-linear complex relationships → Axis Bias

LOF → Local Outlier Factors

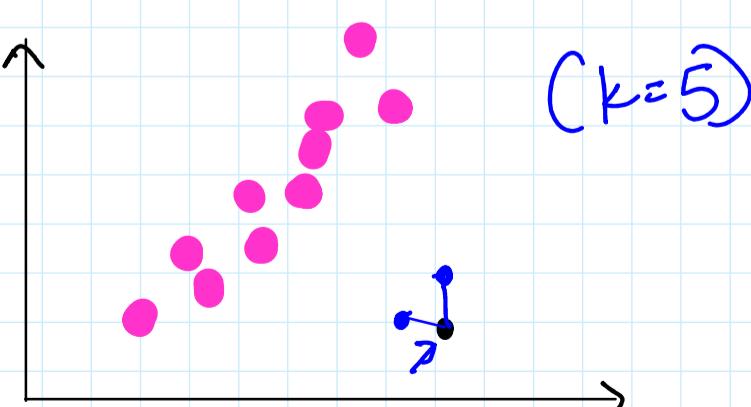


LOF = $\text{LOF} = \frac{\text{Ratio of density of } k\text{-NN}}{\text{to density of observation itself}}$

$\text{LOF} > 1$ $\left[\begin{array}{l} \text{if } k\text{-NN (points)} > \text{f point*} \\ \Rightarrow \text{(Point has possibility of being outlier)} \end{array} \right] \Rightarrow \text{Neighbours are dense}$

< 1 $\left[\begin{array}{l} \text{if } k\text{-NN (points)} < \text{f point*} \\ \Rightarrow \text{that point is not outlier} \end{array} \right]$

Toy Example



$$\text{LOF} = \frac{\text{NR}}{\text{DR}}$$

Density \bullet^5 = area of \circ with 5 nearest observation

* Density \bullet^5 = area of \circ with radius = their 5th nearest neighbour

*
inverted

{ Density = area of O with radius = their 5th nearest neighbour